

```

1 Suggestion with Lucene index
2 -----
3 Shankar:
4     analyzer,
5     tokenizer,
6     ...
7
8
9     >> vector based storage.
10    >> bag of word models.
11
12    document: has words.
13
14    [document as vector of inverted index.]
15
16    d1: t11, t12, ..., t1k1
17    d2: t21, t22, ..., t2k2
18    d3
19
20    dn
21    -----
22    doc/ terms
23         t1  t2  t3  ... tK
24    d1: s1  s2
25    d2:
26
27    where K = union (ki)
28
29    -----
30    score => normalized tf-idf.
31    -----
32    interverted index: to avoid sparse 0 zeros.
33    -----
34
35    >> Tokenizer: return all the token terms for a document(strem of token).
36
37    d1: Shankar has large family with 3 teenage kids.
38
39    stream of ([shankar, has, large, family, with, 3, teenage, kids])
40
41    filters: stop words [a, an, the, has, numeric, ...]
42
43    filter takes input of (tokenStream)
44
45    Analyzer: StandardAnalyzer, WhiteSpace,....
46
47    ngram analyzer: returns stream of ngrams from the document.
48
49    -----
50
51
52
53
54
55

```